

# A Multi-Modality Evaluation of the Reality Gap in Autonomous Driving Systems

Stefano Carlo Lambertenghi  
*Technical University of Munich, fortiss*  
 Munich, Germany  
 stefanocarlo.lambertenghi@tum.de

Mirena Flores Valdez  
*Technical University of Munich*  
 Munich, Germany  
 mirena.flores@tum.de

Andrea Stocco  
*Technical University of Munich, fortiss*  
 Munich, Germany  
 andrea.stocco@tum.de

**Abstract**—Simulation-based testing is a cornerstone of Autonomous Driving System (ADS) development, offering safe and scalable evaluation across diverse driving scenarios. However, discrepancies between simulated and real-world behavior—known as the reality gap—challenge the transferability of test results to deployed systems. In this paper, we present a comprehensive empirical study comparing four representative testing modalities: Software-in-the-Loop (SiL), Vehicle-in-the-Loop (ViL), Mixed-Reality (MR), and full real-world testing. Using a small-scale physical vehicle equipped with real sensors (camera and LiDAR), and its digital twin, we implement each setup and evaluate two ADS architectures (modular and end-to-end) across diverse indoor driving scenarios involving real obstacles, road topologies, and indoor environments. We systematically assess the impact of each testing modality along three dimensions of the reality gap: actuation, perception, and behavioral fidelity. Our results show that while SiL and ViL setups simplify critical aspects of real-world dynamics and sensing, MR testing improves perceptual realism without compromising safety or control. Importantly, we identify the conditions under which failures do not transfer across testing modalities and isolate the underlying dimensions of the gap responsible for these discrepancies. Our findings offer actionable insights into the respective strengths and limitations of each modality and outline a path toward more robust and transferable validation of autonomous driving systems.

**Index Terms**—autonomous driving; reality gap; virtual testing; real-world testing; vehicle-in-the-loop; mixed-reality.

## I. INTRODUCTION

To ensure the safety and reliability of autonomous driving systems (ADS) before deployment in public environments, rigorous system-level testing is indispensable [1]–[3]. A common industrial practice for ADS validation follows a two-phase testing pipeline. First, ADS components—such as perception and planning modules—are trained using real-world driving data and evaluated within virtual environments via simulation-based testing, also known as simulation-in-the-loop (SiL). Subsequently, the ADS undergoes real-world testing on real vehicles on closed tracks up to public roads [4]–[7]. Real-world testing, while more faithful, is costly, time-consuming, and constrained in scope and safety. Despite their scalability, simulations cannot fully replicate real-world physical phenomena, such as sensor noise, actuator delays, and environmental complexity. The resulting mismatch is known as the *reality gap* [1] and hinders the transferability of findings to real-world ADS, undermining their trustworthiness [8].

Various strategies have been proposed to mitigate the reality gap. Some aim to increase simulation fidelity through high-precision modeling (digital twins [9]–[11]), others address specific gap dimensions, such as the perception gap, by translating simulated sensor outputs into more realistic versions using generative models [8], [12]–[15]. However, these methods are limited to model-level testing [16], [17] and do not capture the system-level interactions between perception, planning, and control modules that govern vehicle motion. As a result, they are susceptible to actuation gaps and often miss critical system-level failures [14], [16], [18]. Other strategies involve vehicle-in-the-loop (ViL) and mixed reality (MR) testing [19], by integrating physical components such as ECUs, small-scale robots, or full vehicles-into simulation loops. While ViL provides closed-loop evaluation, it remains partially virtualized and fails to capture real-world imperfections, such as sensor noise and lighting variability (perception gap) [20]. MR partially mitigates this by injecting virtual elements (e.g., obstacles) into real sensor data, enriching scenario realism. Prior system-level studies using small-scale robots investigating failure transferability [8], [21]–[25] have primarily documented the existence of the reality gap, without isolating their root causes or comparing how different test modalities influence gap reduction.

To this aim, in this paper, we conduct an empirical study of the reality gap in autonomous driving by comparing four representative testing modalities: SiL, ViL, MR, and full real-world execution (RW). Our goal is to characterize the dimensions of the reality gap—specifically along actuation, perception, and behavioral fidelity—and to assess the degree to which each testing setup retains ADS behavior relative to real-world ground truth behavior. While prior work has evaluated or mitigated specific aspects of the reality gap, a broader evaluation spanning multiple testing modalities and ADS architectures remains unaddressed.

To investigate the impact of different testing modalities, we implemented a modular ROS-based evaluation framework that supports direct comparisons across synthetic, hybrid, and physical testing conditions. Our setup integrates both modular and end-to-end ADS architectures on a small-scale platform equipped with real sensors (camera and LiDAR) and its digital twin. We conduct hundreds of tests across matched driving scenarios with shared road layouts, obstacle placements, and

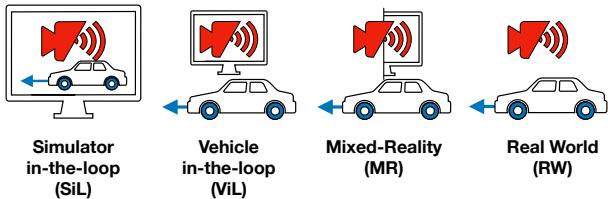


Fig. 1: System-level testing modalities for ADS.

environmental conditions, allowing direct attribution of performance differences to the testing modality.

Our findings reveal that: (i) SiL underestimates real-world variability due to idealized dynamics; (ii) ViL improves actuation realism but retains perception limitations; (iii) MR offers the best perceptual fidelity by blending virtual elements into real sensor data. By isolating the effects of actuation, perception, and behavior on the reality gap, we find that critical failures often manifest differently across configurations, with the perception gap playing a greater role in behavioral divergence than actuation discrepancies. This underscores the importance of testing methods that retain real-world sensor complexity. Our results reveal the limitations of conventional simulation and point to MR as a practical middle ground between fidelity and scalability.

Our paper makes the following contributions:

**Evaluation Framework.** We provide a ROS framework for comparing SiL, ViL, MR, and RW testing for E2E and modular ADS, which is available [26].

**Empirical Study of the Reality Gap.** We present a systematic analysis of the reality gap in ADS testing across behavior, actuation, and perception fidelity, isolating which failures transfer across test modalities. We show that MR testing uniquely replicates real-world system failures, outperforming SiL and ViL across all metrics.

## II. BACKGROUND

### A. Autonomous Driving Systems

Architecturally, ADS can be divided into two categories: end-to-end (E2E) systems and modular systems. E2E systems rely on deep neural networks (DNNs) that directly map camera inputs to driving commands such as steering, throttle, and braking. Once trained, models like NVIDIA’s Dave-2 [27] or InterFuser [28] infer vehicle control actions from raw sensor data without intermediate representations. In contrast, modular ADS architectures such as Pylot [29], Transfuser [30], and Autoware [31] decompose the driving task into distinct components such as perception, planning, and control [32]. The perception module processes raw sensor data (e.g., LiDAR) to detect relevant objects and position them in the perceived environment map. The planning module uses this map to select a safe and feasible route, which the control module executes through low-level actuation commands. As both architectures are actively used and researched [32], [33], we include both in our evaluation to ensure broader applicability of our findings.

### B. Reality Gap Dimensions in ADS Testing

Figure 1 depicts the various system-level testing modalities for ADS. Real-world (RW) testing, conducted on closed-loop tracks or public roads, remains the gold standard for final validation. It exposes the ADS to real-life conditions [4], [6], but is expensive, logically complex, and time-consuming.

To support earlier development stages, simulation-based testing (SiL) offers a scalable and safe environment for experimenting across diverse scenarios. However, simulation introduces a reality gap, a mismatch between simulated and real-world behavior, largely due to limitations in replicating physical sensing and actuation with high fidelity.<sup>1</sup> We refer to the former as the *perception gap*, i.e., the inability of simulated sensors to accurately replicate the real-world sensors. The latter, the *actuation gap*, reflects discrepancies between the vehicle dynamics modeled in simulation and those exhibited by physical vehicles. Together, these issues contribute to the *behavior gap*, where the actions of the ADS in simulation diverge from its behavior in real-world scenarios.

To reduce this gap, vehicle-in-the-loop (ViL) methods embed a real vehicle into a simulated environment, enabling realistic actuation and closed-loop evaluation. While ViL helps address the actuation gap, it typically relies on synthetic sensor inputs and thus remains vulnerable to perception inaccuracies [12], [13]. Mixed reality (MR) testing extends ViL by blending simulated elements directly into real-world sensor streams (e.g., camera images, LiDAR point clouds). This approach preserves the physical characteristics of sensor signals and vehicle dynamics, aiming to jointly mitigate both the perception and actuation gaps.

## III. REALITY GAP EVALUATION FRAMEWORK FOR ADS

To date, a systematic assessment of the various dimensions of the reality gap, or the relative importance of the mitigation strategies, is missing, possibly due to the lack of a standardized evaluation framework. To address this, in this section, we introduce a framework designed to evaluate the transferability of system-level tests across different execution modalities, including SiL, ViL, MR, and real-world (RW) closed-loop testing. Our evaluation targets both E2E and modular ADS configurations for lane-keeping and obstacle-avoidance tasks.

### A. Real-World Setup

1) *Small-Scale Vehicle:* For our RW experiments, we use a small-scale vehicle based on the Donkey Car<sup>TM</sup> open-source framework [35], a widely adopted testbed for ADS research in both simulation and field settings [36]–[38]. The vehicle is equipped with a front-facing 8MP Sony IMX219 RGB camera and a Time-of-Flight (ToF) sensor providing LiDAR-based depth at  $256 \times 192$  resolution and up to 5 m range.

<sup>1</sup>While the term reality gap is also used in the literature to refer to the realism or real-world likelihood of test scenarios in scenario-based testing [34], this aspect is beyond the scope of our work.

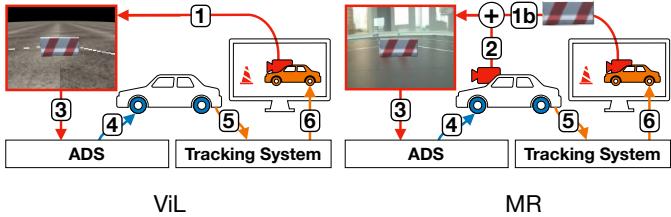


Fig. 2: Data flow and processing steps for ViL and MR.

2) *Testing Tracks*: We conduct experiments in two indoor environments. *Room Nominal* is a dedicated  $6 \times 6$  m robotics lab with minimal background objects, providing a controlled setting. Its  $4 \times 4$  m closed-loop track includes five right and two left curves ( $10^\circ$ – $90^\circ$ ), marked by 10 cm white lane margins and a central dotted line. In contrast, *Room Generalization* is a larger, triangular multi-purpose room ( $20 \times 10 \times 10$  m) with visually complex backgrounds. Its  $6 \times 5$  m stadium-shaped track has two semi-circular curves connected by straights, bounded by narrower 3 cm lane markings and no central line. The floor color also differs, adding perceptual variability.

3) *Tracking system*: The tracking module determines the position of the vehicle and obstacles within the environment using a high-precision motion capture system. It utilizes four Vicon Bonita 10 cameras placed at the corners of each track with an extra margin of 1 m (both for Room Nominal and Generalization), which emit infrared light and detect reflections from retro-reflective markers attached to the vehicle and obstacles. Each object is uniquely identified based on the spatial configuration of its markers at 100 Hz. The Vicon tracker software manages the cameras, computes object poses, and broadcasts them as UDP messages. These messages are received by our framework, which integrates the data into the rest of the system.

#### B. Software setup

1) *Simulator*: The Donkey Car framework provides a high-fidelity Unity3D [39] simulator that models the physical vehicle using the Nvidia PhysX engine [40]. We extended the simulator by developing a new Unity scene that procedurally generates track layouts using Catmull-Rom splines [41], based on real-world lane margin data from Room Nominal and Generalization. To replicate the physical environments, we applied high-resolution (48 MP) floor images as textures and configured the lane markings to match each sandbox. In addition, we modified the simulator to support a virtual depth sensor that mimics the behavior of the ToF system used in our real-world setup. The virtual sensor produces depth data at the same resolution and provides intrinsic matrix parameters consistent with those retrieved by ToF system.

2) *Testing Framework*: To support SiL, ViL, MR, and RW experimentation in a consistent setup, we integrate the simulator and physical platform into a ROS-based software framework [42]. ROS is chosen for its modular, platform-agnostic architecture, which allows the framework to be adapted across different hardware platforms, ADS designs, and

simulation tools without major structural changes. Core functionalities are implemented in dedicated ROS nodes. These nodes communicate via standard topics, enabling modular configuration and straightforward component replacement.

In the SiL modality, both the sensor stream and vehicle dynamics are handled entirely in simulation. The simulator provides perception data to the ADS, and control outputs are executed directly within the game engine, with no involvement of physical hardware or tracking. The framework simply interfaces the simulator with the ADS, routing image data in and receiving control commands back, and receiving telemetry data for modular positioning and experiments monitoring.

In the RW modality, perception data originates exclusively from the physical sensors onboard the vehicle. No simulated data is used. The tracking system is used only to monitor the vehicle’s position and does not interact with the simulator. The framework sends the real sensor data to the ADS and routes the predicted driving commands to the vehicle.

For the ViL and MR modalities, Figure 2 shows how the framework manages data flow between the simulator, the physical vehicle, and the ADS. Simulated data is first generated by the simulator. In the ViL setting, the full simulated sensor stream (e.g., camera or LiDAR) is used as input ①. In the MR setting, only selected features or objects are rendered and extracted ①(b), then merged with real sensor data from the physical vehicle ② to form mixed-reality inputs. The resulting sensor data—fully simulated in ViL or blended in MR—is processed by the ADS ③, which outputs control commands for the physical vehicle ④. The vehicle’s motion is tracked in real time ⑤, and its pose is fed back to the simulator ⑥ to maintain alignment between virtual and real environments. This forms a closed-loop system in which the ADS continuously perceives and acts on synchronized sensor data and vehicle dynamics.

In the rest of this section, we describe the main functionalities necessary to enable these testing procedures.

**Control Modules.** For all testing levels, when using a modular ADS, the waypoint and speed commands produced by the ADS must be converted into actuation primitives. This is handled by two interfaces. The *Waypoint Follower* computes steering and throttle commands to reach target  $(x, y)$  waypoints using a pure pursuit algorithm. A brake command is issued when the target is reached. It outputs commands in the format  $[throttle, steering, brake]$ . The *PID Speed Controller* receives the target, current speed, and the control commands. It modifies the throttle command component to maintain the target speed.

**Simulator Interface.** This component can retrieve rendered sensor outputs such as RGB and depth images. When testing in SiL, it accepts control commands (throttle, steering, brake) and optionally a throttle multiplier to drive the simulated vehicle. Finally, the interface publishes feedback (pose, velocity, obstacle positions) to be used by ADSs that require localization or monitoring behavior.

**Tracking Interface.** This interface continuously tracks the vehicle and any physical obstacles using the external motion capture system. It publishes the vehicle pose, obstacle array,



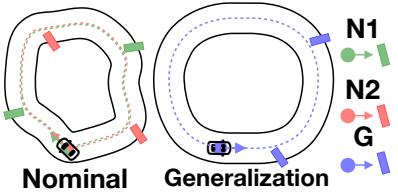


Fig. 4: Testing scenarios used to answer RQ<sub>1</sub>.

straight segments shortly after turns, providing clear visibility and sufficient reaction time. This low complexity setup serves as a baseline for evaluating basic lane keeping and obstacle avoidance. Scenario N2 increases difficulty by placing one obstacle directly on a turn and another immediately after it, within the vehicle’s expected path. This configuration reduces visibility windows and imposes tighter spatial constraints, challenging both perception and planning.

In Room Generalization, for generalization (G), we reuse a printed track featuring turns and straight sections absent from the nominal layout. This new location introduces a perception shift by exposing the camera to unseen background regions. One obstacle is placed on a straight segment (fully visible) and another on a curve (partially occluded). This setup is designed to test the generalization of reality gap mitigation techniques (ViL and MR) on an unseen environment, visual characteristics, and object placements.

### C. Automated Driving Systems

1) *End-to-End*: The end-to-end ADS is a Dave-2 [27] imitation learner trained on real-world camera data to follow the center lane and shift laterally to avoid obstacles with a 0.8m buffer. We recorded 50 laps in Room Nominal at 50 FPS with varying obstacle placements, applied horizontal flipping, and obtained nearly 180k samples. Some laps included recovery from off-track behavior. The training was performed using existing guidelines [27], with a batch size of 64, and a learning rate of 0.0001 for 500 epochs, and the Adam [44] optimizer to minimize the MSE loss, with early stopping (patience 30,  $\Delta$  MSE < 0.05). The final model achieved 0.02 MSE.

2) *Modular*: The modular ADS features a LiDAR-based perception module, a lattice planner, and a control stack with waypoint following and PID speed regulation. This architecture reflects the modular principles adopted in real-world stacks such as Apollo [45] and Autoware [31], albeit simplified for a small-scale testbed. In particular, the perception module adopts a LiDAR DBSCAN clustering approach for obstacle detection [46], a technique widely used in autonomous driving systems [47]–[49]. The planning module uses a Frénet-frame lattice planner, in line with established motion planning methods employed in research and industrial pipelines [50]–[52].

### D. Experimental Setup Validation

For ViL and MR, the real vehicle’s motion must be continuously mapped into the simulation, requiring real-time

execution to keep the simulated state and generated sensor data synchronized with the physical vehicle. Across all domains, the E2E or modular pipeline must execute fast enough to deliver control commands (simulated or real) without latency, as delays could distort the observed behavior. To verify this, we benchmarked the framework over 5000 samples, measuring execution time for every component. All modalities (SiL, ViL, MR, RW) ran in real time: the simulator achieved  $\approx$  80 FPS ( $\approx$  12.5 ms per frame) on our workstation (AMD Ryzen 5, 16 GB RAM, NVIDIA RTX 2060). The tracking system has been measured (at the simulator-level) to provide state updates at  $\approx$  99.93 Hz ( $\approx$  10.01 ms). Since the onboard camera and LiDAR sensors operate at 20 FPS (50 ms per frame), we configure the simulator to run its camera at 20 FPS and to avoid unnecessary system load, we cap the global simulation frame-rate at 60 FPS. The E2E ADS executed in  $\approx$  15 ms on average, and the modular ADS in  $\approx$  20 ms. Inference begins as soon as a new sensor message arrives and can run in parallel if needed, though both pipelines consistently completed before the next frame. Control commands, measured at the vehicle-level, were issued every  $\approx$  50.5 ms, remaining stable across all experiments.

To avoid packet loss, we use a dedicated TP-Link Archer C6 dual-band gigabit WLAN router (867 Mbit/s 5 GHz + 300 Mbit/s 2.4 GHz, 4 Gigabit LAN ports) configured as a local-only network. Vehicle–workstation communication uses TCP over the 5 GHz link (Intel-8265AC, 867 Mbit/s), while tracking–workstation communication is performed via wired Gigabit Ethernet. JPEG-compressed camera images ( $\approx$  800 kB) and compact control commands (< 30 bytes) were observed to fit within each 20 Hz cycle, allowing  $\approx$  50 image transmissions and thousands of control messages without observed packet loss.

We also validated the accuracy of the tracking system using a calibrated reference object with five retro-reflective markers mounted in a fixed, non-coplanar geometry. The employed Vicon system performs tracking by matching such known geometries to objects in space, allowing it to estimate a unique 6-DoF pose (position and three-axis orientation) for each tracked object. To verify accuracy, we measured the distance between two markers of the calibration object with a known spacing of 240 mm, repeating the measurement five times at 30 s intervals in both Room Nominal and Room Generalization. The mean measured distance was 240 mm, with a mean absolute error of 0.40 mm and a standard deviation of 0.0253 mm, confirming sub-millimeter position accuracy. Because orientation is derived from the same geometry matching process, a correct match of the marker configuration also guarantees accurate three-axis orientation. These results demonstrate accurate vehicle localization and synchronization for ViL and MR experiments.

### E. Procedure and Metrics

1) *RQ<sub>1</sub> (behaviour gap)*: To address **RQ<sub>1</sub>**, we evaluate both ADS architectures (E2E and modular) across four domain configurations: the real world (ground truth), SiL (baseline),

ViL, and MR. Experiments are conducted for scenarios N1, N2 (for nominal), and G (for generalization).

In the nominal setting, for each ADS, we run the RW testing modality on scenarios N1 and N2 at least five times and continue until we obtain at least four successful executions. Then, for the SiL, ViL, and MR testing modalities, we run the same number of tests as in RW and quantified the differences in trajectory fidelity, task completion, lane-keeping, and obstacle avoidance. More in detail, we assess system behavior using two categories of metrics. First, we measure trajectory similarity using the Fréchet distance [53] between each run and its real-world counterpart. The Fréchet distance accounts for both spatial proximity and the order of points along the trajectory, making it well-suited for autonomous driving where both path shape and sequence matter. We compute it between each SiL, ViL, or MR run and the corresponding RW run to quantify trajectory deviation. This metric has been widely adopted in prior work [13], [53] to assess behavioral similarity in autonomous systems. Second, we evaluate task performance through several sub-metrics. The *failure rate* is defined as the proportion of runs resulting in either a collision with an obstacle or a lane departure. We also report the absolute number of *obstacle crashes* and *out-of-road* events. To characterize failures, we compute the *completion rate* as the percentage of the track driven before a failure. For SiL, ViL, and MR, the *completion rate* is reported relative to the RW execution. This setting resulted in 92 executions.

In Scenario G, we only evaluate RW, SiL, and MR modalities, omitting ViL, as the room’s features are expected to mainly affect perception. Each configuration is executed three times, and the same set of metrics, trajectory similarity, failure rate, completion rate, and driving quality, are used to assess behavioral consistency across domains. This setting results in 18 executions.

Finally, for the modular ADS, which separates perception and control, we conduct an ablation study. In this configuration, we bypass the perception module and provide ground-truth obstacle locations directly to the planner in both the SiL and RW modalities. This isolates the impact of actuation by removing the perception gap, an analysis not possible for the E2E system, where perception and control are coupled together. By re-purposing the results of the RW and SiL models in scenarios N1 and N2, this setting results in 16 new executions (8 runs x 2 domains).

2) *RQ<sub>2</sub> (actuation gap)*: To address **RQ<sub>2</sub>**, we quantify the actuation gap by comparing the physical response of the real vehicle to that of a simulated vehicle under identical control inputs. We limit the evaluation to SiL, ViL, and RW, omitting MR as this modality does not change actuation fidelity beyond what is already captured by ViL. Each experiment was repeated five times to capture variability, and data were logged at 100 Hz, yielding large frame-level sample sizes. Statistical analysis is applied only to per-frame metrics, where sample sizes support meaningful significance testing. We design five test scenarios spanning both low-level actuation and high-level control, described next. **Forward Motion.** Constant throttle

is applied with steering fixed at zero and braking disabled. Throttle commands are taken directly from the RW-recorded ROS messages with matched timestamps. Trials run for three seconds using throttle values of 0.34, 0.365, and 0.39, which represent the lowest throttle at which the vehicle starts moving, to the highest throttle at which the vehicle can perform the experiment whilst remaining in the tracking area. We compute the total distance traveled, average speed, and trajectory, and compare both speed and trajectory distributions. This setup results in 15 executions (3 speed values x 5 repetitions).

**Steering Motion.** With the throttle fixed at 0.365 and braking disabled, selected based on the medium throttle value identified in *Forward Motion experiments*. Steering commands are issued across six values: {-1.0, -0.6, -0.3, 0.3, 0.6, 1.0}, representing the full steering range, from full left to full right. As with throttle, steering inputs are replayed from recorded ROS messages to ensure the same timing across domains. We compare turning radius via circle fitting, raw trajectories, and trajectory distributions. This setup results in 30 executions (6 steering values x 5 repetitions).

**Braking Motion.** Following forward motion, braking is triggered 35 cm before a 2.0 m goal. Steering remains at zero, and throttle inputs match those from the forward motion test. Braking events are triggered based on each domain’s local pose estimate, so that stopping occurs at the same physical location in RW, SiL, and ViL. We measure braking distance, approach speed, and deceleration. This setup results in 15 executions (3 throttle values x 5 repetitions).

**PID Speed Control.** This scenario evaluates closed-loop speed regulation. The throttle is fixed at 1.0 while a PID controller adjusts output via a throttle multiplier. We replay the exact RW speed requests for each phase to maintain identical target profiles across domains. The vehicle drives in a circular path (steering at 0.6) across four 10-second speed phases: 0.4, 0.8, 0.6, and 0 m/s. We report comparisons of per-phase and overall speed errors. Trajectory similarity is not computed due to the fixed path. This setup results in 5 executions.

**Waypoint Following.** We test tracking performance on pre-defined waypoint paths (throttle fixed at 0.365), ranging from simple (single point) to complex (wide turns, sharp turns, S-shapes). The same actuation module and waypoint-following logic are used in all domains, running closed-loop control on local pose estimates. All waypoint-following tests were performed before adding lane markings in Nominal, as perception is not involved in this experiment. We compute Fréchet distances among trajectories. This setup results in 30 executions (6 paths x 5 repetitions).

3) *RQ<sub>3</sub> (perception validity and gap)*: To evaluate *perception input validity*, we used two test scenarios. All experiments are repeated five times from a fixed initial pose in each domain to capture variability.

**Obstacle Placement.** Static obstacles are placed at 0.4–1.6 m distances, either centrally (single) or symmetrically (dual) within the sensor FoV. For cameras, we compute IoU between manually annotated real and mixed bounding boxes (10 samples/modality). For LiDAR, we calculate the Euclidean







Concerning LiDAR perception, MR also outperforms SiL. For obstacle detection, SiL shows a mean depth error of 0.108 m, a maximum of 1.068 m, and a standard deviation of 0.179 m. MR reduces these to 0.039 m (−64%), 0.713 m (−33%), and 0.090 m, respectively, yielding higher accuracy.

Overall, these results confirm that SiL introduces substantial perception gaps, particularly for camera input, which contribute to behavior failures. MR narrows these gaps, better aligning sensor inputs with real-world distributions and enabling more accurate behavior replication.

**RQ<sub>3</sub>:** *MR significantly reduces perceptual discrepancies compared to SiL, improving camera SSIM from 0.21 to 0.93 and reducing LiDAR error by over 60%, resulting in realistic sensor inputs that closely match real-world data and enable accurate system behavior replication.*

#### D. Threats to Validity

1) *Internal validity:* We compared all ADS under identical parameter settings. One threat to internal validity concerns our custom implementation. However, this was unavoidable as no similar evaluation frameworks are available, to the best of our knowledge. Another threat may be due to our data collection phase and training of ADS, which may exhibit a large number of misbehaviors if trained inadequately or with poor quality data. We mitigated this threat by training and fine-tuning the best publicly available driving models, which performed consistently in nominal RW conditions. The RW-trained E2E model completed scenarios N1 and N2 in all five RW runs without failure (Table 1, Nominal–E2E–Real), indicating adequate training. The LiDAR-based modular pipeline failed in N2 due to late-appearing obstacles, as its perception relies on deterministic clustering rather than learned models; however, the ablation study shows the planning module operates correctly when decoupled from LiDAR inputs. Failures in scenario G stem from domain shift rather than poor training: the E2E ADS succeeds in SiL but fails in RW, which would not occur if the model were fundamentally incapable of solving the task. Overall, these results confirm that our comparisons are not confounded by under-trained or low-quality models.

2) *External validity:* We used a limited number of ADS architecture models in our evaluation, which we mitigated by covering representative ADS architectures. We considered only two physical tracks and a scaled-down platform, which may not capture all real-world physics (e.g., suspension dynamics or vehicle mass distribution). However, our goal is not to model absolute performance but to assess relative fidelity across modalities. Donkey Car was used as a proxy for full-size ADS also in previous studies [8], [14], [21]–[23], [57] and uniquely satisfies our requirements for studying transferability between simulated and real-world testing of ADS. Other platforms, such as DeepPiCar [58] and JetRacer [59], lack integrated simulators; Roboracer [60] offer only low-fidelity physical simulations that do not realistically capture real-world driving dynamics. AWS DeepRacer [61] is tightly integrated

with AWS infrastructure and is primarily designed for reinforcement learning use cases, which are outside the scope of this work. In contrast, Donkey Car has been successfully adopted in numerous real-vehicle autonomous driving studies [8], [14], [21]–[23], [36]–[38], [57], making it a practical and cost-effective experimentation platform.

We acknowledge that the availability of multiple and diverse tracks and obstacle configurations would be desirable. However, our selection of scenarios meant to isolate and evaluate core ADS capabilities, namely lane-keeping and obstacle avoidance, under controlled, repeatable conditions. Hence, generalizability to other physical settings or RC platforms might not hold or may hold partially. We use a scaled-down platform, which may not capture all real-world physics (e.g., suspension dynamics or vehicle mass distribution). However, our goal is not to model absolute performance but to assess relative fidelity across modalities.

3) *Reproducibility:* All software artifacts and results are available in our replication package and appendix [26]. To replicate our study, however, two physical assets are needed, i.e., a Donkey Car and a tracking system.

## VI. DISCUSSION

### A. Dissecting and Addressing the ADS Reality Gap

Our study highlights persistent differences between SiL, ViL, MR, and RW testing, underscoring the need to improve transferability across these environments. While simulation remains indispensable for ADS validation, our results show it is insufficient in isolation. Two major challenges were observed, namely the *perception gap* and the *actuation gap*.

To mitigate the former in SiL, without relying on physical vehicles, recent neural rendering approaches have been proposed. Examples include generative AI for translating simulated into photorealistic images [8], [12]–[15], [62] and diffusion models for generating realistic operational design domains [63]–[65]. Some are already integrated in simulators such as NVIDIA Omniverse [66]. While promising, these methods still suffer from correctness issues (e.g., artifacts, hallucinations), increase runtime cost, and leave the actuation gap and simulator unreliability [67], [68] unaddressed.

On the other hand, the actuation gap is better addressed through hybrid testing. ViL and MR, as evaluated in this work, virtually eliminate actuation mismatches, with MR additionally reducing perception errors. However, these setups are not a substitute for SiL experiments as running tests on real vehicles and hardware remains costly and resource-intensive, even if less than RW. Instead, we view ViL and MR as complementary, highlighting the need to develop strategies to prioritize which scenarios merit RW execution.

Our results highlight meaningful differences across ADS types. In SiL, E2E systems appear under-confident while modular pipelines appear overconfident. Under generalization scenarios that introduce distribution shifts, the trend reverses, with E2E failing in RW while modular ADS proves more

robust. To address this, we suggest a staged strategy: SiL for early validation and coverage, ViL for refining control, and MR for perception fidelity. At the same time, MR setups are more expensive and harder to deploy than SiL, motivating our release of a modular and well-engineered framework.

Finally, our findings must be interpreted with respect to the experimental setup: the vehicle is small-scale, tracks and scenarios are simple, and our open-source simulator, particularly its rendering, is not fully photorealistic. These factors may have amplified the transfer gap, especially for perception-heavy ADS pipelines. Nonetheless, we expect similar issues, though at different magnitudes, even with industrial-grade simulators and full-scale vehicles.

#### B. Implications for testing, debugging, and monitoring

Hybrid testing offers promising opportunities to improve validation, debugging, and monitoring of ADS. For *validation*, we propose the use of ViL and MR to replay critical or failure-inducing cases originally detected in SiL [69], [70]. This approach supports more trustworthy assessments of safety and behavior and enables more representative test generation [71]. From a *debugging* perspective, MR provides a valuable trade-off. Unlike pure simulation, it avoids artificial failures; unlike RW testing, it enables repeatable, cost-effective, and versatile experimentation with physical vehicles. This makes MR particularly suitable for analyzing rare or complex issues such as braking delays or occluded pedestrian responses, although the transferability of such insights to full-scale RW vehicles remains open. For *monitoring*, ViL and MR allow evaluation under realistic latencies, noise, and actuation constraints. Since many state-of-the-art monitoring tools [72]–[75] are assessed only in simulators, hybrid environments help bridge this gap by enabling controlled crash reproduction, evaluation under adverse conditions, and the training of more robust monitoring systems. Early steps in this direction include the works of Ayerdi et al. [76] and Huang et al. [77].

## VII. RELATED WORK

#### A. Reality Gap Assessment Studies

The reality gap has been subjected to active research in many fields, including robotics, automotive, and artificial intelligence. For a comprehensive survey, we refer the reader to Hu et al. [78], while notable contributions are also provided in the software engineering community [1]–[3], [8], [67]. Concerning empirical studies on the reality gap, Stocco et al. [8] compare ADS lane-keeping models in simulated and physical environments, highlighting critical shortcomings that contribute to the gap. In this work, beyond assessing the gap between SiL and RW, we also evaluate mitigation strategies such as ViL and MR. Similarly, Gao et al. [79] propose MultiTest, a physically-aware object insertion framework for testing the robustness of fusion-based perception systems, while Gao et al. [80] outline key challenges in benchmarking AI-enabled multi-sensor fusion across diverse conditions. However, these studies focus primarily on perception robustness within simulated

environments and do not extend to full-system evaluation. In contrast, our work provides, to the best of our knowledge, the first unified, system-level assessment of ViL and MR as reality gap mitigation strategies, explicitly analyzing their impact across perception, actuation, and behavior during live execution in SiL, MR, and RW settings under consistent and controlled conditions. However, the transferability to real-world vehicles is not assessed.

#### B. Reality Gap Mitigation Studies

Concerning solutions to mitigate the gap, researchers have proposed a variety of strategies. One common approach involves the use of digital twins, which aim to replicate real-world vehicle dynamics and sensor characteristics with high fidelity [9]. Alternatively, search-based tuning of simulator parameters can be employed using real-world logs [25]. Another technique is domain randomization, which improves generalization by varying environmental parameters such as lighting, weather, or road conditions during training [20], [81]. However, domain randomization and adversarial training are typically applied to models trained only in simulation. This is not the case for ADS, as they are trained on real-world data.

To address the perception gap, researchers have explored the use of generative AI techniques to translate simulated data into more realistic representations. These efforts include the generation of photorealistic images [8], [12]–[15], realistic LiDAR point clouds [82], and methods for synthesizing operational design domains with high visual fidelity [63], [64], [83], [84]. While existing methods enhance perceptual realism, they typically run offline at the model level or in SiL, targeting single sensor modalities. In contrast, our work evaluates perception fidelity during live, system-level execution across SiL, MR, and RW, directly comparing how perceptual gaps affect the autonomous system system behavior.

## VIII. CONCLUSIONS

This paper presents a comprehensive empirical study of the reality gap in autonomous driving, analyzing Software-in-the-Loop (SiL), Vehicle-in-the-Loop (ViL), Mixed-Reality (MR), and real-world execution across both modular and end-to-end driving systems. Our goal is to isolate the dimensions of the reality gap, namely, behavior, actuation, and perception, and assess how each modality reflects real-world behavior.

Our findings reveal that the reality gap is multifaceted and modality-dependent. SiL often misrepresents system behavior, failing when the real system succeeds, and vice-versa. ViL mitigates actuation errors but leaves perception gaps unresolved. MR, with simulated obstacles and RW perception, is the only testing approach that more consistently captures both perceptual and behavioral fidelity, matching real-world outcomes in both nominal and generalization scenarios.

Our openly available framework and results offer a foundation for the development of next-generation ADS testing and validation solutions, promoting the adoption of hybrid setups that combine SiL, ViL, and MR to enable efficient, scalable ADS evaluation.



